# 现代程序设计第12周作业

谢奕飞 20377077

## 代码

### WriteHead函数

```python
def WriteHead(head,dir_path):
    '''
    写入表头
    '''
    if not os.path.exists(dir_path):
        os.mkdir(dir_path)
    with open(dir_path+'/data.csv','w',encoding='utf8',newline='') as f:
        writer=csv.writer(f)
        writer.writerow(head)
```

### GetPage函数

```python
def GetPage(url,headers):
    '''
    获取页数
    '''
    response=requests.get(url=url,headers=headers)
    soup=BeautifulSoup(response.text,'lxml')
    return int(soup.select('a[class="zpgi"]')[-1].get_text())
```

### Producer函数

```python
def Producer(q:Queue,url,headers):
    '''
    生产者
    '''
    response=requests.get(url=url,headers=headers)
    soup=BeautifulSoup(response.text,'lxml')
    soup_list=soup.select('a[class="tit f-thide s-fc0"]')
    href_list=[]
    for s in soup_list:
        href_list.append('https://music.163.com'+str(s['href']))
    q.put(href_list)
```

### Consumer

```python
def Consumer(q:Queue,headers,dir_path):
```

```
2       '''
3       消费者
4       '''
5       head=['id','title','image','author_id','author','description','coun
    t','play','add','share','comment']
6       urls=q.get()
7       if urls!=None:
8           for url in urls:
9               response=requests.get(url=url,headers=headers)
10              soup=BeautifulSoup(response.text,'lxml')
11              id=url.split('id=')[-1]
12              title=soup.select('.tit')[0].get_text()[1:]
13              image_url=soup.select('img[class="j-img"]')[0]['data-src']
14              image=DownloadImage(image_url,dir_path+'/images')
15              author_id=soup.select('a[class="s-fc7"]')[0]['href'].split('id
    =')[-1]
16              author=soup.select('a[class="s-fc7"]')[0].get_text()
17              description=soup.select('p')[1].get_text()
18              count=soup.select('span[id="playlist-track-count"]')[0].get_te
    xt()
19              play=soup.select('strong[id="play-count"]')[0].get_text()
20              add=soup.select('a[class="u-btni u-btni-fav"]')[0]['data-coun
    t']
21              share=soup.select('a[class="u-btni u-btni-share"]')[0]['data-c
    ount']
22              comment=soup.select('span[id="cnt_comment_count"]')[0].get_tex
    t()
23              with open(dir_path+'/data.csv','a',encoding='utf8',newline='')
    as f:
24                  writer=csv.writer(f)
25                  writer.writerow([id,title,image,author_id,author,descripti
    on,count,play,add,share,comment])
```

## DownloadImage函数
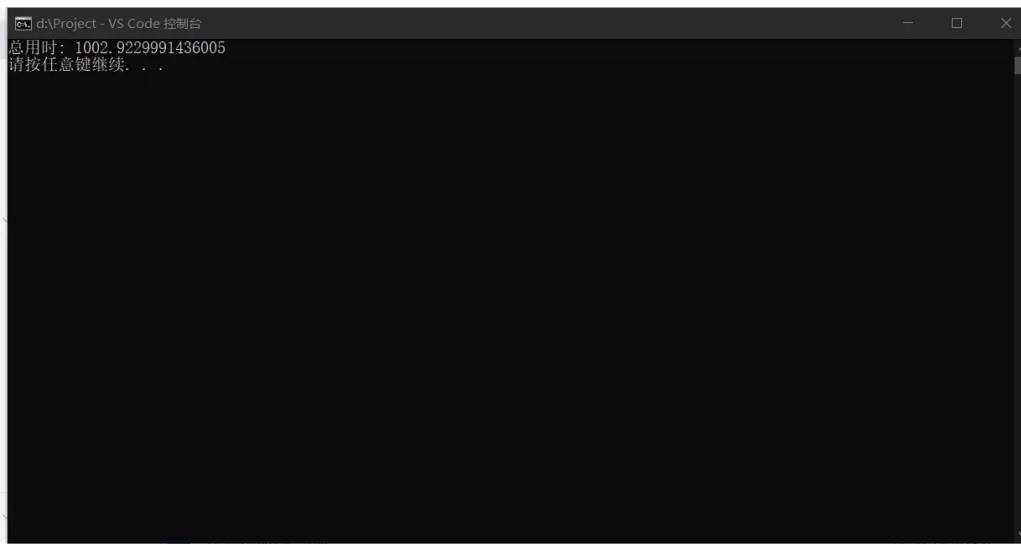
```
1   def DownloadImage(url,dir_path):
2       '''
3       下载图片
4       '''
5       if not os.path.exists(dir_path):
6           os.mkdir(dir_path)
7       name=url.split('/')[-1]
8       img_path=dir_path+'/'+name
9       try:
10          urllib.request.urlretrieve(url,filename=img_path)
11          urllib.request.urlcleanup()
12      except:
13          return 'error'
14      return img_path
```
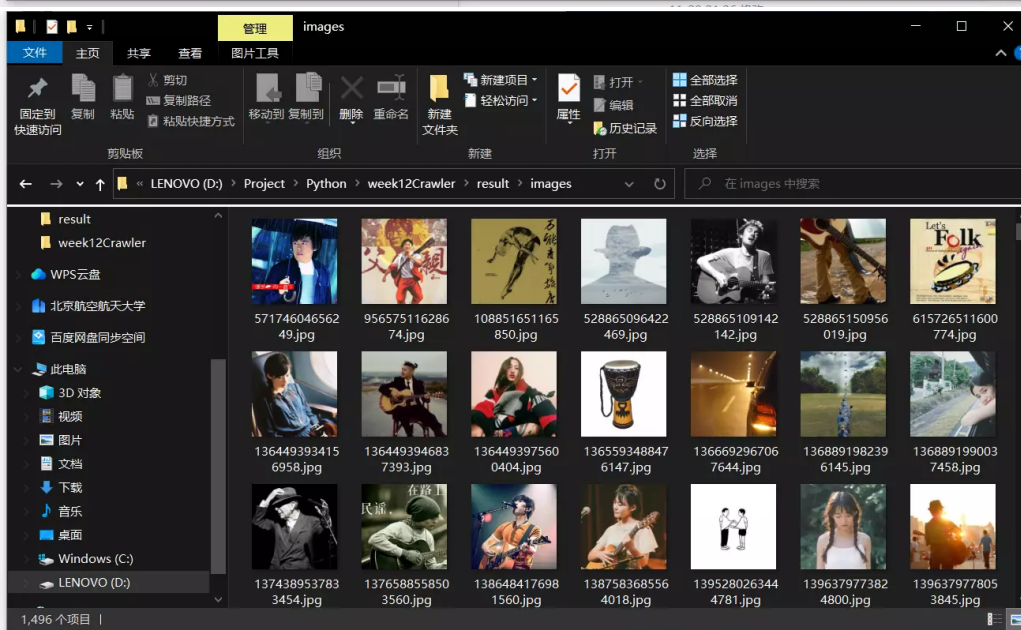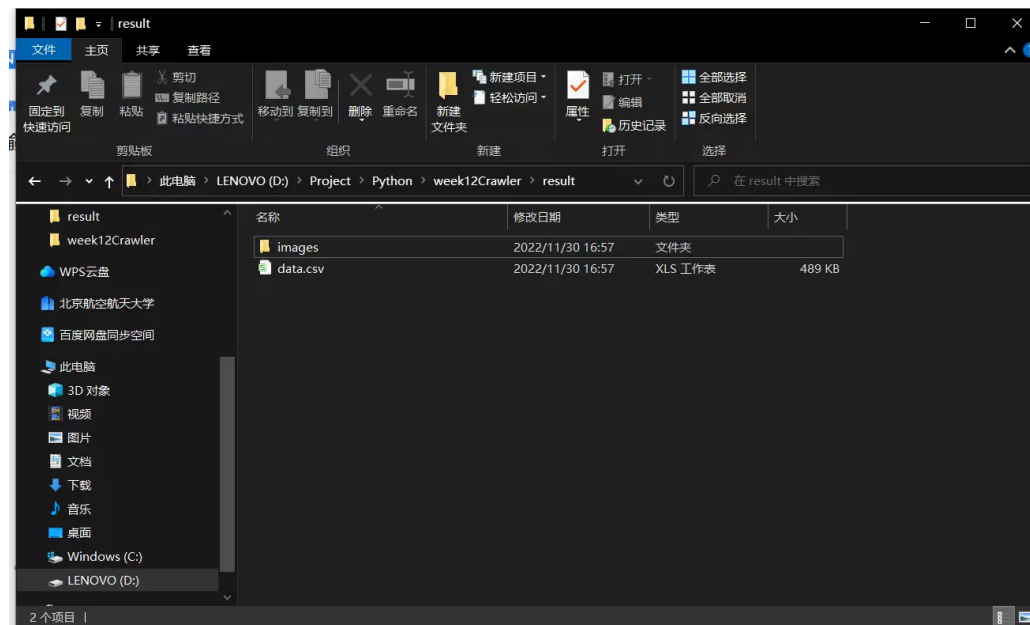
## main函数

```
 1  if __name__=='__main__':
 2      t_start=time.time()
 3      url='https://music.163.com/discover/playlist/?order=hot&cat=%E6%B0%91%
        E8%B0%A3&limit=35&offset=0'
 4      headers={'User-Agent':'Mozilla/5.0 (Windows NT 10.0; Win64; x64) Apple
        WebKit/537.36 (KHTML, like Gecko) Chrome/107.0.0.0 Safari/537.36 Edg/107.
        0.1418.35'}
 5      dir_path='D:/Project/Python/week12Crawler/result'
 6      head=['id','title','image','author_id','author','description','coun
        t','play','add','share','comment']
 7      WriteHead(head,dir_path)#写入表头
 8      N=GetPage(url,headers)#获取页数
 9      urls=[]
10      for i in range(N):
11          urls.append(f'https://music.163.com/discover/playlist/?order=hot&c
        at=%E6%B0%91%E8%B0%A3&limit=35&offset={i*35}')
12      q=Queue()
13      #Producer(q,url,headers)
14      #Consumer(q,headers,dir_path)
15
16      #爬
17      plist,clist=[],[]
18      for url in urls:
19          p=Thread(target=Producer,args=(q,url,headers,))
20          plist.append(p)
21      for i in range(N):
22          c=Thread(target=Consumer,args=(q,headers,dir_path,))
23          clist.append(c)
24
25      for p in plist:
26          p.start()
27      for c in clist:
28          c.start()
29      for p in plist:
30          p.join()
31      for c in clist:
32          q.put(None)#主进程发信号结束，但要给每一个consumer准备
33      for c in clist:
34          c.join()
35      t_finish=time.time()
36      print('总用时: {}'.format(t_finish-t_start))
```

## 运行结果

总用时: 1002.9229991436005
请按任意键继续. . .

用时1002.92s

# 附录–完整代码

```python
import requests
from bs4 import BeautifulSoup
import time
from queue import Queue
from threading import Thread
import csv
import os
import urllib.request

def WriteHead(head,dir_path):
    '''
    写入表头
    '''
    if not os.path.exists(dir_path):
        os.mkdir(dir_path)
    with open(dir_path+'/data.csv','w',encoding='utf8',newline='') as f:
        writer=csv.writer(f)
        writer.writerow(head)

def GetPage(url,headers):
    '''
    获取页数
    '''
    response=requests.get(url=url,headers=headers)
    soup=BeautifulSoup(response.text,'lxml')
    return int(soup.select('a[class="zpgi"]')[-1].get_text())

def Producer(q:Queue,url,headers):
    '''
    生产者
    '''
    response=requests.get(url=url,headers=headers)
    soup=BeautifulSoup(response.text,'lxml')
    soup_list=soup.select('a[class="tit f-thide s-fc0"]')
    href_list=[]
    for s in soup_list:
        href_list.append('https://music.163.com'+str(s['href']))
```

```python
38          q.put(href_list)
39
40  def Consumer(q:Queue,headers,dir_path):
41      '''
42      消费者
43      '''
44      head=['id','title','image','author_id','author','description','count','play','add','share','comment']
45      urls=q.get()
46      if urls!=None:
47          for url in urls:
48              response=requests.get(url=url,headers=headers)
49              soup=BeautifulSoup(response.text,'lxml')
50              id=url.split('id=')[-1]
51              title=soup.select('.tit')[0].get_text()[1:]
52              image_url=soup.select('img[class="j-img"]')[0]['data-src']
53              image=DownloadImage(image_url,dir_path+'/images')
54              author_id=soup.select('a[class="s-fc7"]')[0]['href'].split('id=')[-1]
55              author=soup.select('a[class="s-fc7"]')[0].get_text()
56              description=soup.select('p')[1].get_text()
57              count=soup.select('span[id="playlist-track-count"]')[0].get_text()
58              play=soup.select('strong[id="play-count"]')[0].get_text()
59              add=soup.select('a[class="u-btni u-btni-fav"]')[0]['data-count']
60              share=soup.select('a[class="u-btni u-btni-share"]')[0]['data-count']
61              comment=soup.select('span[id="cnt_comment_count"]')[0].get_text()
62              with open(dir_path+'/data.csv','a',encoding='utf8',newline='') as f:
63                  writer=csv.writer(f)
64                  writer.writerow([id,title,image,author_id,author,description,count,play,add,share,comment])
65
66  def DownloadImage(url,dir_path):
67      '''
68      下载图片
69      '''
70      if not os.path.exists(dir_path):
71          os.mkdir(dir_path)
72      name=url.split('/')[-1]
73      img_path=dir_path+'/'+name
74      try:
75          urllib.request.urlretrieve(url,filename=img_path)
76          urllib.request.urlcleanup()
77      except:
78          return 'error'
79      return img_path
80
```

```python
if __name__=='__main__':
    t_start=time.time()
    url='https://music.163.com/discover/playlist/?order=hot&cat=%E6%B0%91%
E8%B0%A3&limit=35&offset=0'
    headers={'User-Agent':'Mozilla/5.0 (Windows NT 10.0; Win64; x64) Apple
WebKit/537.36 (KHTML, like Gecko) Chrome/107.0.0.0 Safari/537.36 Edg/107.
0.1418.35'}
    dir_path='D:/Project/Python/week12Crawler/result'
    head=['id','title','image','author_id','author','description','coun
t','play','add','share','comment']
    WriteHead(head,dir_path)#写入表头
    N=GetPage(url,headers)#获取页数
    urls=[]
    for i in range(N):
        urls.append(f'https://music.163.com/discover/playlist/?order=hot&c
at=%E6%B0%91%E8%B0%A3&limit=35&offset={i*35}')
    q=Queue()
    #Producer(q,url,headers)
    #Consumer(q,headers,dir_path)

    #爬
    plist,clist=[],[]
    for url in urls:
        p=Thread(target=Producer,args=(q,url,headers,))
        plist.append(p)
    for i in range(N):
        c=Thread(target=Consumer,args=(q,headers,dir_path,))
        clist.append(c)

    for p in plist:
        p.start()
    for c in clist:
        c.start()
    for p in plist:
        p.join()
    for c in clist:
        q.put(None)#主进程发信号结束，但要给每一个consumer准备
    for c in clist:
        c.join()
    t_finish=time.time()
    print('总用时：{}'.format(t_finish-t_start))
```